**Original Research**

**Title:** Predicting postoperative risks using large language models

Bing Xue[1,2†], Charles Alba[2†], Joanna Abraham[2,3], Thomas Kannampallil[2,3], Chenyang Lu[1,2,3*]

[1]McKelvey School of Engineering, Washington University in St. Louis, St. Louis, 63130, Missouri, USA, [2]AI for Health Institute, Washington University in St. Louis, St. Louis, 63130, Missouri, USA. [3]School of Medicine, Washington University in St. Louis, St. Louis, 63130, Missouri, USA

[†]These authors contributed equally to this work

[*]Corresponding author: E-mail(s): lu@wustl.edu

**Abstract**

**Background:** Predicting postoperative risk can inform effective care management and planning. We explored large language models (LLMs) in predicting postoperative risk through clinical texts using various tuning strategies.

**Methods:** Records spanning 84,875 patients from Barnes Jewish Hospital (BJH) between 2018 and 2021, with a mean duration of follow-up based on the length of postoperative ICU stay less than 7 days, were utilized. Methods were replicated on the MIMIC-III dataset. Outcomes included 30-day mortality, pulmonary embolism (PE) and pneumonia. Three domain adaptation and finetuning strategies were implemented for three LLMs (BioGPT, ClinicalBERT and BioClinicalBERT): self-supervised objectives; incorporating labels with semi-supervised fine-tuning; and foundational modelling through multi-task learning. Model performance was compared using the area under the receiver operating characteristic curve (AUROC) and the area under the precision recall curve (AUPRC) for classification tasks and mean squared error (MSE) and $R^2$ for regression tasks.

**Results**: Cohort had a mean age of 56.9 (sd: 16.8) years; 50.3% male; 74% White. Pre-trained LLMs outperformed traditional word embeddings, with absolute maximal gains of 38.3% for AUROC and 14% for AUPRC. Adapting models through self-supervised finetuning further improved performance by 3.2% for AUROC and 1.5% for AUPRC Incorporating labels into the finetuning procedure further boosted performances, with semi-supervised finetuning improving by 1.8% for AUROC and 2% for AUPRC and foundational modelling improving by 3.6% for AUROC and 2.6% for AUPRC compared to self-supervised finetuning.

**Conclusions:** Pre-trained clinical LLMs offer opportunities for postoperative risk predictions with unseen data, and further improvements from finetuning suggests benefits in adapting pre-trained models to note-specific perioperative use cases. Incorporating labels can further boost performance. The superior performance of foundational models suggests the potential of task-agnostic learning towards the generalizable LLMs in perioperative care. To facilitate the convenient deployment of our models, we have made them publicly available on [HuggingFace](HuggingFace).

**Introduction**

Perioperative care is of paramount importance in improving healthcare quality and patient safety in surgical settings. More than 10% of surgical patients experience major postoperative complications, including infections, such as pneumonia, and blood clots, such as pulmonary embolism (PE) or deep vein thrombosis (DVT).[1-7] These can lead to increased mortality and need for intensive care, extended hospital stays, and higher healthcare costs.[8] Many of these complications are preventable,[1,9-11] highlighting the importance of early identification of patient risk factors. Considering that over half of all healthcare adverse events are due to medical management rather than underlying diseases,[9-11] the development of predictive models for perioperative care has the

potential to support early intervention for improved outcomes. To date, machine learning tools developed for perioperative care utilize mainly numerical and categorical variables or time-series measurements from Electronic Health Records (EHR).[1-4]

Clinical text, which contain vital patient information and details concerning scheduled procedures, remains under-utilized in predicting complications and outcomes in perioperative care.[9] These notes play a crucial role in the decision-making process that impacts the course of a patient's perioperative journey,[12] including the preparation for surgery, the transfer of patients to operating rooms, and the prioritization of clinicians' tasks,[12,13] henceforth underlining their function towards achieving safe patient outcomes.[9,12]

The emergence of LLMs has paralleled the rise of foundational models (FMs), which empower technologies like ChatGPT.[14-18] In the context of healthcare, FMs can perform many different tasks after being pre-trained on large datasets. [14,17-22] Despite these promising developments, using LLMs for clinical prediction tasks presents unique challenges. First, most pre-trained models are trained on limited clinical corpora compared to other domains due to the low-resource setting of clinical notes, which are not readily open-sourced.[22-24] Second, clinical notes can contain inconsistent descriptions of varying lengths and are strongly abbreviated with specialized medical terminologies[25] not frequented amongst ordinary text used for training of pre-trained large language models (LLMs). The latter has prompted researchers to develop clinical LLMs using publicly available EHR databases[19,20,26] or biomedical texts.[21]

Recognizing the continuing need for predictive models in perioperative care and the vital role of clinical notes in this field, we explored the practical application of LLMs to predict postoperative risks. This includes comparing the performance of pre-trained LLMs with traditional word embeddings and experimenting across various finetuning

strategies, including the development of FMs, which may enable clinicians to utilize these models in a task-agnostic manner.[27]

**Methods**

*Settings and Data Sources*

Our main dataset was sourced from electronic anaesthesia records (Epic) for all adult patients undergoing surgery at the Barnes Jewish Hospital (BJH) spanning 4 years (2018 to 2021). The dataset's size was 84,875 patient records. In terms of the textual characteristics, the clinical notes contained a vocabulary size of 3,203, with a mean word and vocabulary lengths of 8.9 (sd: 6.9) and 7.3 (sd: 4.4), respectively. The clinical notes are smart text records with descriptions of scheduled procedures derived from anaesthetic records. The non-textual characteristics of patients are detailed in Appendix A1.1. The internal review board of Washington University School of Medicine in St. Louis (IRB#201903026) approved the study with a waiver of patient consent. We replicated the models and techniques on MIMIC-III,[20] a publicly available dataset of critical care patients at the Beth Israel Deaconess Medical Center between 2001 and 2012. To truly replicate the approach and settings employed in BJH's clinical notes, we utilized the long-form descriptive texts of procedural codes in MIMIC-III. More details, including the characteristics of MIMIC-III's clinical notes, are detailed in Appendix A1. This manuscript follows the Transparent Reporting of a Multivariable Prediction Model for Individual Prognosis or Diagnosis (TRIPOD) reporting guideline. A comparison of characteristics of the cohorts from these datasets can be found in Appendix A1.

*Process of data collection*

BJH notes were derived pre-operatively from smart text records containing descriptions and details pertaining to scheduled procedures within each patient's anaesthetic records,

which were pulled for our study. Outcomes were determined based on a combination of laboratory values, dialysis events, structured anaesthesia assessments, billing data, nurse flow-sheets, and ICD-10 diagnosis codes. More details on the data collection process for our dataset and the replication on MIMIC-III could be found in Appendix A2.

### *Data Preprocessing*

The clinical text was provided by BJH in a de-identified and censored format before being handed to the authors for analysis as a method to preserve patient privacy, which is important since our models publicly available on Huggingface. This process involved removing sensitive or uniquely identifiable information, including textual descriptions of conditions or procedures that could reveal patient identities, and reformatting the text of scheduled procedures to include only common tokens in an arbitrary order. The pre-processing steps are detailed in Appendix A3.

### *Outcome Variables*

Our main outcomes included 30-day mortality, pulmonary embolism (PE), and pneumonia, in addition to deep vein thrombosis (DVT), delirium, and acute knee injury (AKI) reported in the appendix in the interest of space. These six outcomes are pertinent to perioperative care, particularly during OR-ICU handoff[1,28]. For the MIMIC-III replication, outcomes included ICU in-hospital mortality, length-of-stay (LoS), 30-day mortality, and 12-hour discharge status, chosen based on previous studies.[2,26,29,31]

### *Models*

We employed BERT and GPT-based LLMs for postoperative prediction, namely clinicalBERT[19], bioClinicalBERT[20], and bioGPT[21]. The architectures and pre-training corpora for each of these models are elaborated in Appendix A4. These models have been tested across representative NLP benchmarks in the medical domains

including Question-Answering tasks benchmarked by the PubMedQA,[32] MedQA[33] and ClinicalQA,[34] recognizing entities from texts,[35] or extracting relations and patient characteristics through text.[36] In this work, we evaluate its performance in predicting postoperative outcomes in the context of perioperative care by examining its utilization in different settings, such as using the pre-trained model without any finetuning or employing a variety of finetuning strategies. We will detail these strategies in the subsequent sections. Parameters used across these different finetuning strategies are detailed in Appendix A4.4. In addition to the clinical LLMs, we included traditional NLP methods as baselines for comparison, including word2vec's continuous bag-of-words (CBOW),[37] doc2vec,[38] GloVe,[39] and FastText.[40]

### Predictors

To ensure a consistent comparison across different methods, including traditional NLP word embeddings, pre-trained LLMs, and their distinct fine-tuned variants, we standardized the approach by using the eXtreme Gradient Boosting Tree (XGBoost)[41] as the predictor for all outcomes. This choice allows us to accommodate a diverse range of input types while leveraging XGBoost's widespread use in healthcare due to its robust performance in various clinical prediction tasks.[1,2] Among the models of the best performing finetuning strategy, we experimented with various predictors, including Random Forest, Logistic Regression, and the feed-forward auxiliary network found in models that incorporate labels into the finetuning process. The range of parameters employed for each predictor are detailed in Appendix A4.5.

### Evaluation metrics and validation strategies

Experiments were stratified into 5 folds for cross-validation. A nested cross-validation approach was used. This means that within each fold, 20% of the data in the

inner loop was designated as a validation set to monitor and select the best hyperparameters within each fold.

For classification tasks, we calculated the area under the receiver operating characteristic curve (AUROC) and the area under the precision-recall curve (AUPRC) to get a comprehensive evaluation of the models' overall prediction performance in the face of class imbalance. For the best-performing models, we also reported their accuracy, sensitivity, specificity, precision, and F-scores in the Appendix A5. For regression tasks, we calculated the mean squared error (MSE) and $R^2$.

Four different experiments were conducted to investigate: (1) how pre-trained clinical LLMs perform on unseen data/tasks, (2) how different finetuning approaches affect perioperative predictions, (3) how our proposed foundational strategy improves model performance, and (4) how various ML predictors perform on textual embeddings.

### Pre-trained LLMs on unseen data/tasks

We first evaluated the general applicability of clinical LLMs in perioperative predictions without developing a bespoke model for each postoperative outcome, as illustrated in Figure 1a.

### Transfer learning: self-supervised finetuning vs semi-supervised finetuning

In the setting of transfer learning of pre-trained LLMs, we explored two approaches: self-supervised finetuning and semi-supervised finetuning. The self-supervised finetuning leverages the information contained within the source domain and exploits it to align the distributions of source and target data. In this approach, we adapted the pre-trained LLMs with our training data in accordance with their existing self-supervised training objectives. For BioGPT, this entails the language modeling task. For ClinicalBERT and BioClinicalBERT, it entails the masked language modeling

(MLM) and Next Sentence Prediction (NSP) objectives. These are elaborated in the Appendix A4.
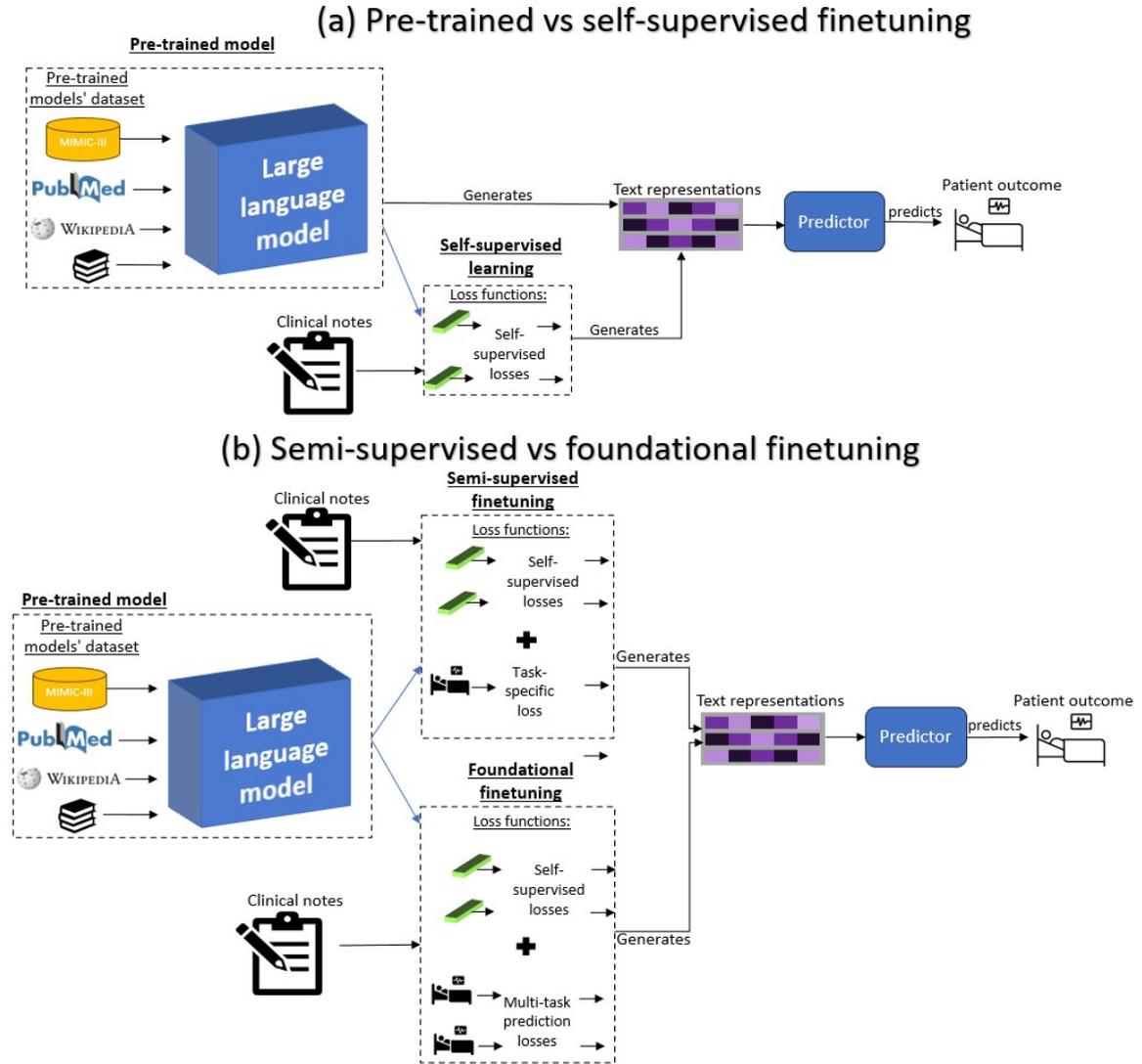


Figure 1. An illustration of the architectures encompassing different tuning strategies experimented in our study. Fig 1a (top) illustrates how using the pre-trained model alone differs from self-supervised finetuning when clinical texts are provided to the pre-trained LLM to refine the model weights with respect to its objective loss functions. Fig 1b (below) illustrates two separate finetuning strategies: semi-supervised finetuning – creating a model that is finetuned under the supervision of a specific outcome; and foundational finetuning – creating a foundational model that is finetuned through a multi-task learning (MTL) objective using all available postoperative labels in the dataset.

Studies have shown that incorporating labels into the finetuning process can lead to improved predictive performances[19,42]. Extending from such studies, we adapt

a semi-supervised finetuning approach in lieu of expected improvements with self-supervised finetuning. This approach further utilizes the labels (postoperative outcomes) of texts – information that was not available during pre-training and not used in self-supervised finetuning. In addition to the self-supervised training objectives mentioned above, an auxiliary fully connected neural network predictor was introduced to supervise the textual embedding to better align with the training labels. Details of the auxiliary network's architecture are detailed in Appendix A4. A $\lambda$ parameter was introduced serving as the coefficient in controlling for the magnitude of the loss between the self-supervised tasks and the newly introduced supervision auxiliary network. The parameters employed under each finetuning strategy, including the optimal choices of $\lambda$ for each task and dataset, are reported in Appendix A4.4.

*Foundational finetuning strategy*

To build a foundational model with knowledge across all tasks, we extended the above-mentioned semi-supervised strategy and propose a task-agnostic finetuning strategy, inspired from Aghajanyan et al.[27] It exploits all possible labels available within the dataset, including but not limited to selected tasks. This makes the finetuning *foundational* in the sense that it solves various tasks simultaneously by employing a multi-task learning framework for knowledge sharing across all available labels in the dataset. Each label is assigned a task-specific auxiliary network wherein the losses across all labels are pooled together. The parameters employed under the foundational finetuning strategy are reported in Appendix A4.4.

*Examining the effects of ML predictors on predictive performance*

After exporting the textual embeddings from LLMs, the choice of machine learning predictor is flexible. As detailed in the 'predictors' sub-section, we investigated how

different predictors influence the predictive performance of the best-performing models.

**Results**

Amongst our dataset, a total of 84,875 patient (mean [SD] age, 56.9 [16.8] years; 50.3% male; 74% White) notes, 1,694 of these cases resulted in 30-day mortality (positive event rate, 2%), 287 had Pulmonary Embolism (PE) (positive event rate, 0.3%), and 475 had pneumonia (positive event rate, 0.6%). Details of the additional outcomes could be found in Appendix A5.

### Pre-trained LLMs on unseen data/tasks

*Table 1. A comparison of traditional NLP models (top) vs pre-trained LLMs (bottom). The results are presented as the mean and 95% confidence interval across all 5-folds. In the interest of space, the results for the additional three experimented outcomes are reported in Appendix A5.1. The best baseline models are underlined, and the best models are **bolded**. As shown amongst the results, the pre-trained LLMs consistently outperform traditional word embeddings.*

| Model | 30-day mortality | | PE | | Pneumonia | |
|---|---|---|---|---|---|---|
| | AUROC | AUPRC | AUROC | AUPRC | AUROC | AUPRC |
| cbow | 0.528 (0.409, 0.648) | 0.023 (0.015, 0.031) | 0.506 (0.418, 0.593) | 0.004 (0.002, 0.006) | 0.526 (0.384, 0.668) | 0.009 (0.001, 0.016) |
| Doc2vec | 0.479 (0.348, 0.611) | 0.021 (0.012, 0.03) | 0.517 (0.466, 0.567) | 0.004 (0.004, 0.004) | 0.495 (0.347, 0.643) | 0.006 (0.003, 0.01) |
| fastText | 0.725 (0.67, 0.781) | 0.05 (0.04, 0.06) | 0.652 (0.602, 0.701) | 0.007 (0.005, 0.01) | 0.696 (0.643, 0.749) | 0.016 (0.008, 0.024) |
| GloVe | <u>0.818</u> <u>(0.807, 0.83)</u> | <u>0.128</u> <u>(0.118, 0.139)</u> | <u>0.664</u> <u>(0.628, 0.701)</u> | <u>0.01</u> <u>(0.007, 0.013)</u> | <u>0.765</u> <u>(0.732, 0.799)</u> | <u>0.04</u> <u>(0.017, 0.063)</u> |
| bioClinicalBERT | 0.85 (0.84, 0.861) | 0.156 (0.138, 0.173) | 0.683 (0.621, 0.745) | 0.008 (0.006, 0.011) | 0.809 (0.785, 0.833) | 0.043 (0.027, 0.059) |
| bioGPT | **0.862** **(0.851, 0.872)** | **0.161** **(0.141, 0.182)** | 0.711 (0.679, 0.743) | 0.011 (0.005, 0.017) | **0.818** **(0.8, 0.837)** | **0.047** **(0.037, 0.058)** |
| ClinicalBERT | 0.855 (0.842, 0.867) | 0.155 (0.137, 0.173) | **0.717** **(0.691, 0.743)** | **0.013** **(0.009, 0.017)** | 0.806 (0.784, 0.827) | 0.04 (0.024, 0.056) |

Table 1 highlights that the use of pre-trained LLMs consistently and significantly surpassed the baseline NLP models. We observed absolute increases that ranged from up to 21.1% in PE to 38.3% for mortality in AUROC. Similarly, increases in the AUPRC ranged from 0.9% in PE to an impressive 14% in 30-day mortality. Results for the additional outcomes can be found in A5.1. A similar magnitude of

performance improvements was noted on the MIMIC-III dataset, as highlighted in Appendix A5.2. This leap in performance highlights the sheer power of pre-trained LLMs in grasping clinically relevant context, even in a setting adjacent 'zero-shot' scenario, when compared to word-based embeddings. Meanwhile, it is worth noting that there was no single pre-trained clinical LLM that always outperformed the rest.

*Transfer learning: self-supervised finetuning vs semi-supervised finetuning vs foundational finetuning*

After finetuning the pre-trained LLMs with perioperative data, we observed absolute improvements ranging from up to 0.6% for 30-day mortality to 3.2% in PE for AUROC. For AUPRC, increases ranged from up to 0.3% in PE to 1.5% in 30-day mortality. These findings were also observed in experiments with the additional outcomes and on the MIMIC-III dataset (Appendix A5), emphasizing the effect of finetuning the weights of these pre-trained models to tailor them to specific clinical domains.

The introduction of labels in the semi-supervised and foundational approaches further enhanced prediction performance compared to the self-supervised method. In contrast to the self-supervised approach which adjusts weights based solely on training texts, the semi-supervised method 'infuses' label information during the finetuning process. As a result, predictive performance is boosted relative to the self-supervised approach since both textual and labelled data were utilized during training. Specifically, when comparing semi-supervised finetuning with self-supervised finetuning, we observed moderate improvements for AUROC ranging from up to 0.7% in 30-day mortality to 1.8% in PE, and for AUPRC ranging from 0.5% in pneumonia to 2% in 30-day mortality, relative to the self-supervised approach. Similar findings were also

observed in the experiments with the additional outcomes, reported in Appendix A5.1, and replicated on MIMIC-III, reported in Appendix A5.2.
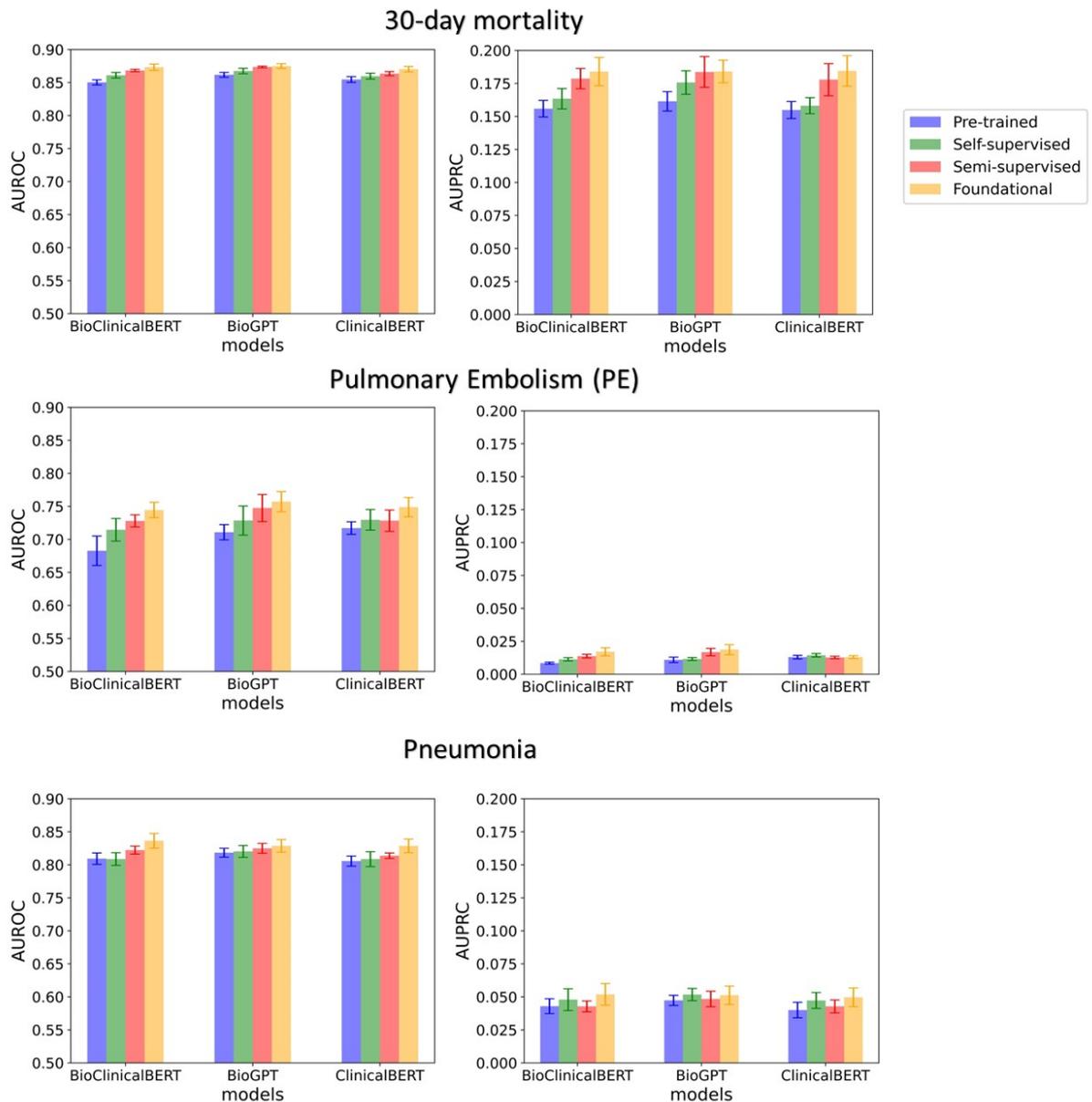


*Figure 2. Comparison of the predictive performance across various models and their respective tuning strategies. Pre-trained LLM consistently outperformed baseline word embeddings. Finetuning improved prediction performance, with the incorporation of labels in the semi-supervised strategy further boosting prediction performances. The model performs best with the foundational strategy, wherein the model was finetuned across all outcomes. (In the interest of space, the figures for the three additional outcomes can be referenced in Appendix A5.1.)*

For foundational finetuning, we incorporated all available labels simultaneously into the finetuning as supervision losses. For our dataset, this encompassed 30-day mortality, PE and pneumonia, acute knee injury, deep vein thrombosis, delirium, acute kidney injury, deep vein thrombosis and delirium. On MIMIC-III, we utilized in-hospital mortality, LOS, 30-day hospitality, and 12-hour discharge status. With the foundational finetuning, we were able to use a single model to outperform various bespoke models for different tasks, whilst observing further increases in AUROC and AUPRCs. For AUROC, the improvement ranged from up to 1.2% in 30-day mortality to 3.6% in PE, when compared with the self-supervised approach. Similarly, AUPRC values increased, ranging from up to 0.4% in pneumonia to 2.6% in 30-day mortality, in comparison to the self-supervised approach. A comparable scale of improvement was observed in the experiments involving the additional outcomes, illustrated in Appendix A5.1, and mirrored on MIMIC-III, reported in Appendix A5.2.

The performance comparison between clinical LLMs varies across tasks. The foundational BioGPT model achieved the best performance for 30-day mortality and PE, whereas the foundational BioClinicalBERT model performed the best for pneumonia, on both AUROC and AUPRC evaluation metrics, as shown in Figure 2. In our MIMIC-III replication, the foundational ClinicalBERT model was the top performer for length-of-stay and 30-day mortality in both AUROC and AUPRC and MSE and $R^2$, respectively. The semi-supervised BioClinicalBERT variant was the best performing model in MIMIC-III's in-hospital mortality in both AUROC and AUPRC, whilst the semi-supervised bioGPT and foundational bioGPT performed best for 12-hour discharge status for AUROC and AUPRC, respectively.

The best performing models are open sourced in Huggingface. With safety being of paramount importance in LLMs, we ran several prompts prior to open sourcing

our model to ensure our model could be safely relied upon during clinical deployment. These prompts can be found in Appendix A6.
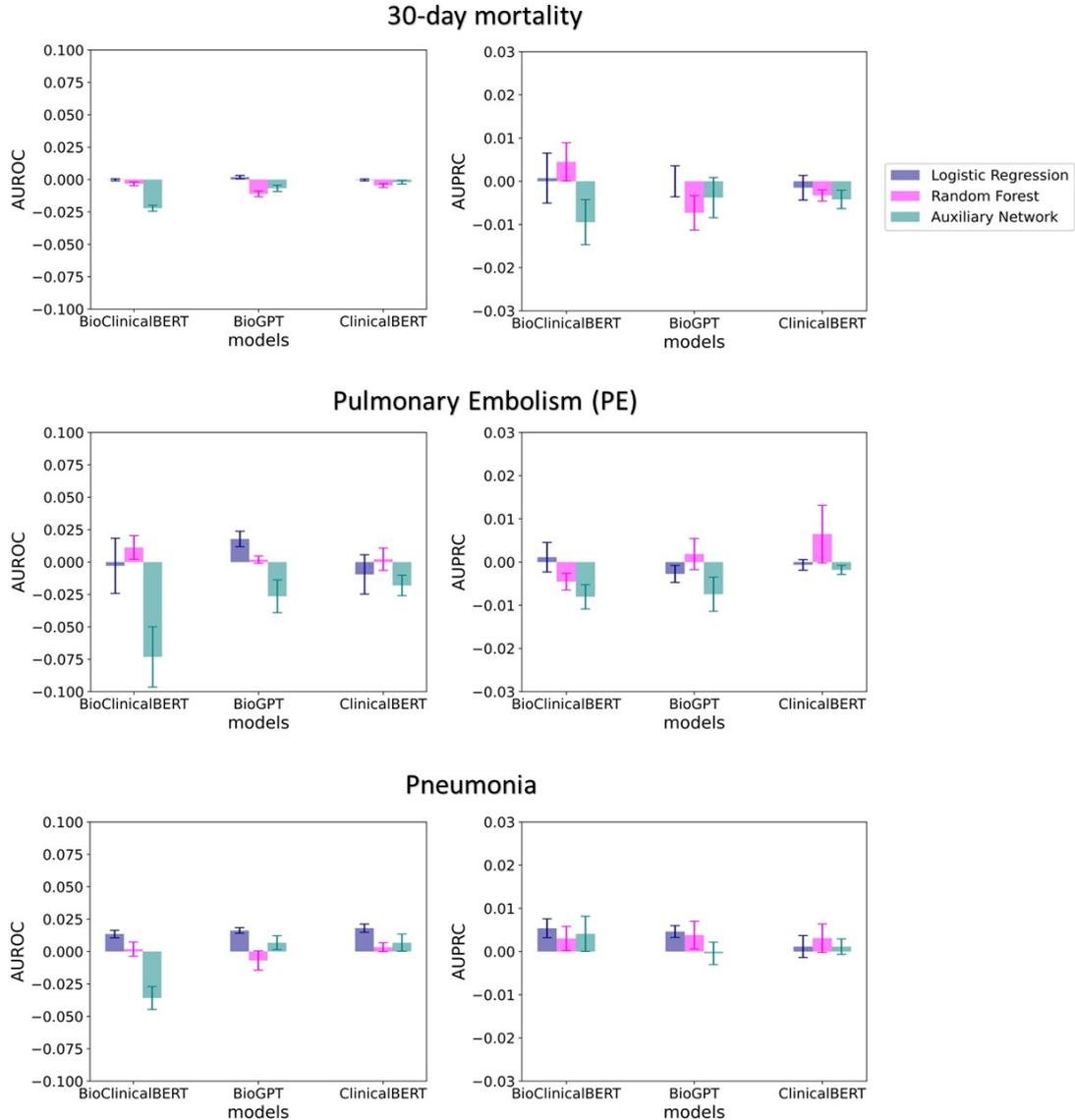


*Figure 3. Comparison of different machine learning classifiers with that of our default XGBoost predictor applied to our textual representations ($\Delta model_{i,j} - XGBoost_i$ with outcome i and model j), including the use of the trained auxiliary layer directly from our foundational model. No single classifier dominated the others across all outcomes and metrics. Surprisingly, the logistic regression classifier performed slightly better than the others, demonstrating that well-tuned language models can generate precise contextual representations to suit a simple classifier. The results the additional outcomes could be referenced in Appendix A5.1.*

To investigate the optimal machine learning model towards textual embeddings, our experiments demonstrated that no single predictor dominantly outperformed the others. Surprisingly, logistic regression performed slightly better than the others, indicating that well-tuned language models can generate precise contextual

14

representations in a linear fashion to suit a simple classifier with better robustness than complex classifiers.

**Discussion**

Improvements in patient safety are of utmost priority in perioperative care, where notes serve as a vital source in clinical decision-making. Yet, Flemons et al[43] noted that a lack of knowledge and the absence of a routine for clinical notes, combined with competing time demands, are the main unintentional factors leading to missed recommendations among clinicians in surgical care. In this context, LLMs could step in and support clinical decisions in perioperative care. While prior AI models utilized tabular EHR data in perioperative care,[1-4] it remains an open challenge to exploit preoperative texts in predicting postoperative outcomes. Our study represents an important step forward in incorporating LLMs into perioperative decision-making, thereby supporting surgical care management, perioperative optimization, and early risk detection.[4]

Our results demonstrate that: (A) pre-trained LLMs can lead to significant improvements over traditional word embeddings, signalling a new era in leveraging LLMs for perioperative care and underscoring their capability adjacent to a 'zero-shot'[24] classification in cases where labelled data may be scarce; (B) finetuning these LLMs can result in further improvements, attributable to the adaptation of these pre-trained models to clinical texts, which may vary in abbreviations and terminology across task-specific datasets; (C) incorporating labels into the finetuning process can improve performance relative to self-supervised finetuning, as tokens linked to specific outcomes are more effectively optimized; and (D) a foundational finetuning strategy yielded the best results, suggesting that a single, comprehensively fine-tuned model

could be effectively deployed to predict a wide range of postoperative outcomes from preoperative clinical notes.

The superior performance of the foundational approach offers tangible benefits in perioperative care. A single model can be fine-tuned to predict multiple outcomes, effectively saving time and resources. Furthermore, foundational models are more generalizable as they have been optimized through knowledge sharing. This alleviates concerns of overfitting to task-specific samples — a core concern in practical applications of AI in medicine.[1-4]

Variations in performance emerged when comparing the results of different machine learning predictors trained on extracted textual embeddings, with no single classifier consistently outperforming the rest. This suggests that conveniently using the final auxiliary layer already in our foundational models may be sufficient to achieve competitive results without the need for training a separate classifier using the LLM's embeddings. This holistic approach facilitates efficient model finetuning and deployment in practice.

Despite demonstrating the potential of LLMs in perioperative care, recent findings of race-based bias in healthcare applications of LLMs require us to acknowledge and use our models and techniques with caution due to potential biases in our data.[44] Firstly, our text data reflect the scheduled procedures of patients. Since these procedures are determined by clinicians who are not immune to personal biases,[45] it is likely that such biases will be replicated in the trained LLMs. Secondly, our BJH data contain a disproportionately low percentage of Hispanic individuals (per Appendix 1), at less than 2%, which could lead the models to learn and replicate patterns predominant in the majority group, potentially reinforcing existing disparities. Hence, by replicating our methods and models on the MIMIC-III dataset, we aim to alleviate these concerns

evolving the association between predictive performance and potential biases unique to our BJH dataset.

Our study also possesses other limitations beyond the bias in the data. First, the quality of the textual data may possess limitations, ranging from incompleteness to potentially missing data, which could potentially impact the performance of our finetuned LLMs. Second, non-textual variables were not used in the models. Adding those variables could potentially improve the model performance.[1] Third, it remains unclear if the observations can be generalized to LLMs in other applications or at various scales. Finally, the impact of the predictions using LLMs on clinical decisions and outcomes need to be studied in a clinical setting. To address these limitations, our ongoing work involves data triangulation across the administrative data, clinical text, and other data to align with high-quality manual health record review provided by National Surgical Quality Improvement Program adjudicators and exploring studies that leverage the LLMs in the context of OR-ICU handoffs.

**Notes**

**References**

1. Xue B, Li D, Lu C, et al. Use of Machine Learning to Develop and Evaluate Models Using Preoperative and Intraoperative Data to Identify Risks of Postoperative Complications. *JAMA Netw Open*. 2021;4(3):e212240. doi:10.1001/jamanetworkopen.2021.2240.

2. Xue B, et al. Perioperative predictions with interpretable latent representation. In: *Proceedings of the 28th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*. 2022;4268-4278. doi:10.1145/3534678.3539190

3. Xue B, et al. Assisting clinical decisions for scarcely available treatment via disentangled latent representation. In: *Proceedings of the 28th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*. 2023;5360-5371. doi: 10.1145/3580305.3599774

4. Xue B, et al. Multi-horizon predictive models for guiding extracorporeal resource allocation in critically ill COVID-19 patients. *J Am Med Inform Assoc*. 2023;30:656-667. doi: 10.1093/jamia/ocac256

5. Hamel MB, Henderson WG, Khuri SF, Daley J. Surgical outcomes for patients aged 80 and older: morbidity and mortality from major noncardiac surgery. J Am Geriatr Soc. 2005;53(3):424-429. doi:10.1111/j.1532-5415.2005.53159.x

6. Turrentine FE, Wang H, Simpson VB, Jones RS. Surgical risk factors, morbidity, and mortality in elderly patients. *J Am Coll Surg*. 2006;203(6):865-877. doi:10.1016/j.jamcollsurg.2006.08.026

7. Healey MA, Shackford SR, Osler TM, Rogers FB, Burns E. Complications in surgical patients. *Arch Surg*. 2002;137(5):611-617. doi:10.1001/archsurg.137.5.611

8. Tevis SE, Kennedy GD. Postoperative complications and implications on patient-centered outcomes. *J Surg Res*. 2013;181(1):106-113. doi:10.1016/j.jss.2013.01.032

9. FitzHenry F, Murff HJ, Matheny ME, et al. Exploring the frontier of electronic health record surveillance: the case of postoperative complications. *Med Care*. 2013;51(6):509-516. doi:10.1097/MLR.0b013e31828d1210

10. Downing NL, Bates DW, Longhurst CA. Physician burnout in the electronic health record era: are we ignoring the real cause? Ann Intern Med. 2018;169(1):50-51. doi: 10.7326/M18-0139

11. Rule A, Bedrick S, Chiang MF, Hribar MR. Length and Redundancy of Outpatient Progress Notes Across a Decade at an Academic Medical Center. JAMA Netw Open. 2021;4(7):e2115334. doi:10.1001/jamanetworkopen.2021.15334

12. Braaf S, Manias E, Riley R. The role of documents and documentation in communication failure across the perioperative pathway. A literature review. *Int J Nurs Stud*. 2011;48(8):1024-1038. doi: 10.1016/j.ijnurstu.2011.05.009

13. Riley R, Manias E. Governing the operating room list. In: The discourse of hospital communication: Tracing complexities in contemporary health care organizations. *London: Palgrave Macmillan UK*; 2007:67-89.

14. Mikolov T, Chen K, Corrado G, Dean J. Efficient estimation of word representations in vector space. *arXiv*. 2013;1301.3781. Preprint posted online Jan 16, 2013

15. Bojanowski P, Grave E, Joulin A, Mikolov T. Enriching word vectors with subword information. *Trans Assoc Comput Linguist*. 2017;5:135-146. doi: 10.1162/tacl_a_00051

16. Wornow M, et al. The shaky foundations of large language models and foundation models for electronic health records. *NPJ Digit Med*. 2023;6:135. doi: 10.1038/s41746-023-00879-8

17. Eriksen AV, Möller S, Ryg J. Use of GPT-4 to diagnose complex clinical cases. *NEJM AI*. 2023;1(1). doi:10.1056/AIp2300031.

18. Lee J, et al. Biobert: a pre-trained biomedical language representation model for biomedical text mining. *Bioinformatics*. 2020;36:1234-1240. doi:10.1093/bioinformatics/btz682

19. Huang K, Altosaar J, Ranganath R. Clinicalbert: Modeling clinical notes and predicting hospital readmission. In *CHIL 2020 Workshop*. 2020.

20. Alsentzer E, Murphy J, Boag W, Weng WH, Jindi D, Naumann T, McDermott M. Publicly available clinical BERT embeddings. In: *Proceedings of the 2nd Clinical Natural Language Processing Workshop*. 2019;72-78. doi:10.18653/v1/W19-1909

21. Luo R, et al. Biogpt: generative pre-trained transformer for biomedical text generation and mining. *Brief Bioinform*. 2022;23:bbac409. doi:10.1093/bib/bbac409

22. Brown T, et al. Language models are few-shot learners. *Adv Neural Inf Process Syst*. 2020;33:1877-1901.

23. Wu H, et al. A survey on clinical natural language processing in the United Kingdom from 2007 to 2022. *NPJ Digit Med*. 2022;5:186. doi: 10.1038/s41746-022-00730-6

24. Zakka C, et al. Almanac—Retrieval-Augmented Language Models for Clinical Medicine. *NEJM AI*. 2024;1(2). doi:10.1056/AIoa2300068.

25. Kuhn IF. Abbreviations and acronyms in healthcare: when shorter isn't sweeter. *Pediatr Nurs*. 2007;33.

26. Johnson AE, et al. Mimic-iii, a freely accessible critical care database. *Sci Data*. 2016;3:1-9. doi:10.1038/sdata.2016.35

27. Aghajanyan A, et al. Muppet: Massive multi-task representations with prefinetuning. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, 2021:5799–5811 doi: 10.18653/v1/2021.emnlp-main.468

28. Fritz BA, et al. Deep-learning model for predicting 30-day postoperative mortality. *Br J Anaesth*. 2019;123:688-695. doi:10.1016/j.bja.2019.07.025

29. Bellini V, Valente M, Bertorelli G, et al. Machine learning in perioperative medicine: a systematic review. *J Anesth Analg Crit Care*. 2022;2(2). doi:10.1186/s44158-022-00033-y.

30. Struja T, et al. Evaluating equitable care in the ICU: Creating a causal inference framework to assess the impact of life-sustaining interventions across racial and ethnic groups. *medRxiv*. Preprint posted online October 13, 2023

31. Luo M, Chen Y, Cheng Y, Li N, Qing H. Association between hematocrit and the 30-day mortality of patients with sepsis: a retrospective analysis based on the large-scale clinical database MIMIC-IV. *PLoS One*. 2022;17(3):e0265758. doi:10.1371/journal.pone.0265758.

32. Jin Q, Dhingra B, Liu Z, Cohen WW, Lu X. PubMedQA: A dataset for biomedical research question answering. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*. 2019;2567–2577. doi:10.18653/v1/D19-1259

33. Yim J, et al. Predicting conversion to wet age-related macular degeneration using deep learning. *Nat Med*. 2020;26:892-899. doi:10.1038/s41591-020-0867-7

34. Mirza FN, et al. Using ChatGPT to facilitate truly informed medical consent. *NEJM AI*. 2024;1(2). doi:10.1056/AIcs2300145

35. Zhu R, Tu X, Huang JX. Utilizing BERT for biomedical and clinical text mining. In: *Data analytics in biomedical engineering and healthcare*. Academic Press; 2021. p. 73-103. doi:10.1016/B978-0-12-819314-3.00005-7

36. Yang X, Yu Z, Guo Y, Bian J, Wu Y. Clinical relation extraction using transformer-based models. *arXiv*. Preprint posted online Jul 19, 2021

37. Church KW. Word2vec. *Nat Lang Eng*. 2017;23:155-162. doi: 10.1017/S1351324916000334

38. Pennington J, Socher R, Manning CD. Glove: Global vectors for word representation. In *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)*. 2014;1532-1543. doi:10.3115/v1/D14-1162

39. Le Q, Mikolov T. Distributed representations of sentences and documents. *PMLR*. 2014;1188-1196.

40. Joulin A, et al. Fasttext.zip: Compressing text classification models. *arXiv*. Preprint posted online Dec 12, 2016.

41. Chen T, Guestrin C. XGBoost: A scalable tree boosting system. In: *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*. 2016:785-794. doi:10.1145/2939672.2939785

42. Chen X, Beaver I, Freeman C. Fine-tuning language models for semi-supervised text mining. In: Proceedings of the 2020 IEEE International Conference on Big Data (Big Data). IEEE; 2020. doi: 10.1109/BigData50022.2020.9377810

43. Flemons K, Bosch M, Coakeley S, et al. Barriers and facilitators of following perioperative internal medicine recommendations by surgical teams: a sequential,

explanatory mixed-methods study. Perioper Med. 2022;11:2. doi:10.1186/s13741-021-00236-x.

44. Logé C, Ross E, Dadey DYA, Jain S, Saporta A, Ng AY, Rajpurkar P. Q-Pain: A Question Answering Dataset to Measure Social Bias in Pain Management. In: Proceedings of the Thirty-fifth Conference on Neural Information Processing Systems (NeurIPS) Datasets and Benchmarks Track (Round 1); June 2021.

45. Tripathi S, Fritz BA, Avidan MS, Chen Y, King CR. Algorithmic bias in machine learning-based delirium prediction. ML4H Extended Abstract Collection. Machine Learning for Health (ML4H); 2022.